

04 NOVEMBER 1999

62



Europäisches  
Patentamt

European  
Patent Office

Office européen  
des brevets

GB 99/3081

REC'D 17 NOV 1999

WIPO PCT

Bescheinigung

Certificate

Attestation

Die angehefteten Unterla-  
gen stimmen mit der  
ursprünglich eingereichten  
Fassung der auf dem näch-  
sten Blatt bezeichneten  
europäischen Patentanmel-  
dung überein.

The attached documents  
are exact copies of the  
European patent application  
described on the following  
page, as originally filed.

Les documents fixés à  
cette attestation sont  
conformes à la version  
initialement déposée de  
la demande de brevet  
européen spécifiée à la  
page suivante.

Patentanmeldung Nr. Patent application No. Demande de brevet n°

98307434.5

**PRIORITY  
DOCUMENT**

SUBMITTED OR TRANSMITTED IN  
COMPLIANCE WITH RULE 17.1(a) OR (b)

Der Präsident des Europäischen Patentamts,  
im Auftrag

For the President of the European Patent Office

Le Président de l'Office européen des brevets  
p.o.

*Aslette Fiedler*

A. Fiedler

DEN HAAG, DEN  
THE HAGUE,  
LA HAYE, LE

09/09/99

This Page Blank (uspto)



Europäisches  
Patentamt

European  
Patent Office

Office européen  
des brevets

**Blatt 2 der Bescheinigung**  
**Sheet 2 of the certificate**  
**Page 2 de l'attestation**

Anmeldung Nr.:  
Application no.:  
Demande n°: 98307434.5

Anmeldetag:  
Date of filing:  
Date de dépôt: 14/09/98

Anmelder:  
Applicant(s):  
Demandeur(s):  
Aston University  
Birmingham B4 7ET  
UNITED KINGDOM  
Amersham Pharmacia Biotech UK Limited  
Little Chalfont, Buckinghamshire HP7 9NA

UNITED KINGDOM  
Bezeichnung der Erfindung:  
Title of the invention:  
Titre de l'invention:  
Gene and protein libraries and production method thereof

In Anspruch genommene Priorität(en) / Priority(ies) claimed / Priorité(s) revendiquée(s)

Staat:  
State:  
Pays:

Tag:  
Date:  
Date:

Aktenzeichen:  
File no.  
Numéro de dépôt:

Internationale Patentklassifikation:  
International Patent classification:  
Classification internationale des brevets:

C12N15/10, C12N15/12, C07K14/47, C12Q1/68, G01N33/68

Am Anmeldetag benannte Vertragsstaaten:  
Contracting states designated at date of filing: AT/BE/CH/CY/DE/DK/ES/FI/FR/GB/GR/IE/IT/LI/LU/MC/NL/PT/SE  
Etats contractants désignés lors du dépôt:

Bemerkungen:  
Remarks:  
Remarques:

The original title of the application reads as follows:  
Gene and protein libraries and methods

*This Page Blank (uspto)*

## GENE AND PROTEIN LIBRARIES AND METHODS

### 5 Introduction

Naturally occurring proteins are capable of specific binding interactions with other proteins and other molecules. It is well known that such proteins can be used as scaffolds and specific amino acid residues changed in order to improve binding properties. The changes required can  
10 be determined by combinatorial chemistry means. The subject is reviewed by Per-Ake Nygren and Mathias Uhlen in Curr. Opin. Struct. Biol. (1997) 7, 463-469, who list cyclic peptides, immunoglobulin-like scaffolds, bacterial receptors, DNA-binding proteins and protease inhibitors as examples of protein scaffolds. The authors conclude that, starting from a suitable  
15 protein domain, the use of a combinatorial approach coupled with powerful selection or screening strategies can be used to obtain novel proteins capable of binding a desired target molecule. But the selection or screening strategies can be difficult. It is this problem that is addressed by the present invention.

20 Zinc fingers are examples of protein scaffolds of the kind described. Zinc fingers are protein motifs ("mini-domains") which interact with double-stranded DNA (some also bind RNA). This interaction is dependent on DNA sequence, thus the interaction is termed to be sequence-specific. The interaction between the zinc finger and its target  
25 DNA sequence is modular: one zinc finger recognises three bases of DNA. Basic rules concerning the interaction were determined early on by structural studies (both X-ray crystallography and NMR spectroscopy) of zinc finger-DNA complexes. In essence, three residues (amino acids) within the zinc finger make base-specific contacts with the DNA. These  
30 three residues differ greatly between different zinc fingers, allowing a limited repertoire of different DNA sequences to be recognised. Early

- 2 -

mutagenesis experiments determined that if these variable residues are changed, a different DNA sequence may be recognised. (A fourth residue sometimes contributes to DNA recognition, but this residue is well-conserved between different zinc finger proteins). In practice then, the zinc finger may be viewed as a molecular scaffold, which orientates the three variable residues suitably to enable them to make base-specific contacts with the DNA.

It would be most advantageous to have available a zinc finger to bind each trinucleotide (3 bases) of dsDNA. Initial attempts to achieve this goal centred on the structure-based design of novel zinc finger proteins. Since 1994 however, several groups have employed combinatorial libraries of zinc finger proteins and/or target DNA sequences to identify novel zinc fingers which bind to the required DNA sequences

One such technique has been developed by Choo and Klug and is described in WO 96/06166 and in PNAS, 91, 11163-11167 and 11168-11172 (1994). A single library of zinc finger genes was constructed. The library was based on a naturally occurring zinc finger protein, Zif 268, which contains three zinc fingers. Only the central finger was randomised at seven positions. The library of genes was cloned as a fusion to the fd phage gene pIII. When expressed, a library of bacteriophage resulted, in which each bacteriophage displayed a randomised zinc finger protein on its surface. In a first stage assay, this library was incubated with a target DNA molecule, and individual clones that bound to the target were purified and sequenced. In a second stage assay, each of those clones selected was incubated with a variety of related DNA sequences in order to further investigate its binding properties. The technique is subject to some inherent disadvantages:

- Deconvolution is not addressed – purification is inherent in the method. The assay results in a pool of a bacteriophage. For identification purposes, each member of that pool must be cultured independently and its DNA sequenced.

- 3 -

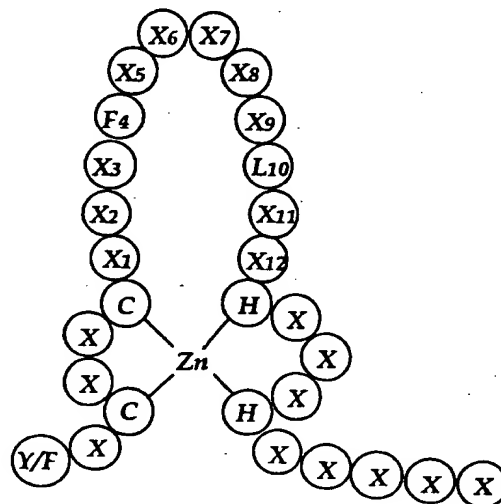
• The experimental end point is determined empirically. While the assay is in progress, it is impossible to determine the number of different phage binding to the target DNA. The end point is therefore determined empirically e.g. by 15 washes. Any zinc finger which binds to the target DNA with sufficient strength to withstand these washes is selected, and a pool of zinc fingers results. There is no in-built mechanism to determine relative binding strengths of zinc fingers within this selected pool; hence the need for a second stage assay.

• Library size. Constructing a library of the size required is technically difficult – indeed, the authors largest library is 200 times smaller than that theoretically required. When expressed therefore, several zinc finger proteins may be omitted.

The present invention addresses these shortcomings.

Zinc fingers are small protein motifs. They form parts of larger proteins, but perform their specific function within those proteins. Zinc fingers exist in tandem arrays: proteins containing between 2 and 37 different zinc fingers have been identified.

In two dimensions, a single zinc finger appears as follows:



20

In this diagram, each circle represents a single amino acid

- 4 -

residue.

The zinc finger is so stable that its structure is unaffected by the replacement of virtually all residues marked "X" with alanine (Michael *et al*, PNAS 89, 4796-4800, 1992). Spaced correctly (as above) the following requirements are all that are necessary for the formation of a zinc finger:

- The 2 cysteine (C) residues
- The 2 histidine (H) residues
- The zinc ion (Zn), which is co-ordinated (bound) by the C and H residues
- Three hydrophobic residues: tyrosine/phenylalanine (Y/F); phenylalanine (F4); leucine (L10).

Zinc fingers bind to nucleic acids - either DNA or RNA. In nature, zinc fingers usually form part of transcription factors, but in the laboratory, it is possible to work with them independently from the rest of these proteins. The zinc finger exemplified herein binds to double-stranded DNA. One zinc finger binds to three bases of DNA (a trinucleotide).

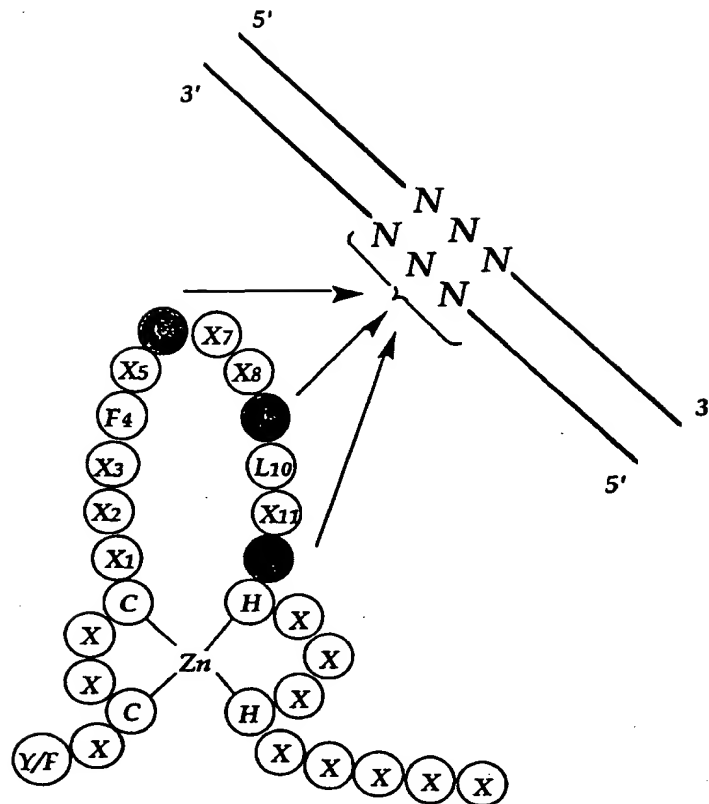
Several zinc fingers are usually linked in tandem. Most frequently, three zinc fingers interact with successive trinucleotides, which means that altogether, the three zinc fingers will interact with (recognise) a specific 9 base pair (bp) sequence of DNA. Each zinc finger will recognise a specific trinucleotide. However, nature has only provided a limited repertoire of zinc fingers, so the number of 9 base pair sequences which can be recognised is very limited.

The mechanism of DNA recognition is sequence-specific and surprisingly simple. Three residues (amino acids) within the zinc finger make contacts (hydrogen bonds or Van de Waal's interactions, for example) with three bases of DNA. Most of these contacts are with one strand of the DNA.

30



- 5 -



Many experiments have shown that if the three interacting residues (here named  $\alpha$ ,  $\beta$  and  $\gamma$ ) are changed, the resulting zinc finger will recognise a different sequence of DNA. Moreover, if a library of zinc finger proteins is made in which  $\alpha$ ,  $\beta$  and  $\gamma$  are randomised, new zinc finger proteins may be identified by screening the library with a specific sequence of DNA.

There are 64 possible trinucleotides:

10

$$\begin{array}{c} \text{Number of trinucleotides NNN} = 4 \times 4 \times 4 = \underline{\underline{64}} \\ | \\ \text{(A,C,G or T)} \end{array}$$

15

Therefore, 64 different zinc finger proteins, each of which binds optimally to one trinucleotide would represent: a complete zinc finger

- 6 -

code. A problem (addressed by this invention) is to develop such a code.

This invention involves applying the principles of combinatorial chemistry to the problem. The key to any combinatorial system (whether biological, chemical or any other system) is deconvolution:  
5 the identification of an active substituent from within a mixture. The key to discovering an optimal zinc finger for each trinucleotide is to identify the optimum combinations of residues  $\alpha$ ,  $\beta$  and  $\gamma$ . There will be an optimum combination of  $\alpha$ ,  $\beta$  and  $\gamma$  for each trinucleotide. By using multiple libraries of zinc fingers, with highly controlled overlap between the libraries,  
10 deconvolution can be achieved without purification.

### The Invention

In one aspect the invention provides a set of libraries of genes which code for proteins which are capable of specific binding  
15 interactions by virtue of amino acid residues at two or more determined positions including a first determined position and one or more other determined positions, which set of libraries consists of:

a) 6 to 20 libraries in which each library has a triplet that codes for one or several but less than 20 amino acids at the said first determined  
20 position, and is randomised at the triplet or triplets coding for the said one or more other determined positions, the arrangement being such that interactions of the proteins coded for by the said 6 to 20 libraries with a specific binding partner identifies a triplet that codes for an amino acid at the said first determined position that takes part in the specific binding  
25 interaction, and

b) 6 to 20 libraries of corresponding design for each of the said one or more other determined positions.

In another aspect the invention provides a set of libraries of proteins, which proteins are capable of specific binding interactions by  
30 virtue of amino acid residues at two or more determined positions including a first determined position and one or more other determined positions,

- 7 -

which set of libraries consists of:

- a) 6 to 20 libraries in which each library has one or several but less than 20 amino acid residues at the said first determined position and is randomised at the said one or more other determined positions, the arrangement being such that interaction of the 6 to 20 libraries with a specific binding partner identifies an amino acid residue at the said first determined position that takes part in the specific binding interaction, and
- b) 6 to 20 libraries of corresponding design for each of the said one or more other determined positions.

In another aspect the invention provides a method of identifying a protein which interacts with a specific binding partner, which method comprises providing a set of libraries of proteins as defined, incubating the specific binding partner with each library of the set, observing specific binding interactions with certain libraries of the set, and using the observations to identify a protein which interacts with the specific binding partner. Preferably, as discussed in more detail below, this method is performed using scintillation proximity assay (SPA) technology.

A library of compounds (e.g. genes or proteins) consists of a plurality of compounds which are all different but which have some characteristic in common. The compounds of the library may be presented either separate or together, in solution or solid phase. In a set of libraries, the compounds of any one library have some characteristic in common but which differentiates them from the compound of each other library of the set.

A specific binding interaction of a protein with another molecule (the specific binding partner) is an interaction mediated by a specified amino acid residue at one or more usually several positions in the protein molecule. The specific binding partner is usually though not necessarily a polymeric molecule, e.g. a nucleic acid (DNA or RNA) or another protein.

In relation to proteins, the statement that a library is

- 8 -

randomised at a determined position is herein used to mean that the library contains a random mixture of all or almost all possible amino acid residues. We say "almost all" because there might be a special reason for omitting one residue e.g. Cys, or a few amino acid residues. In relation to genes,  
5 the statement that a triplet is randomised is herein used to indicate a triplet NNN (where N is any nucleotide) or a triplet that is capable of coding for all or almost all the amino acids.

The term protein is herein used to encompass any chain of two or more amino acid residues.

10 The term polynucleotide is herein used to encompass any chain of three or more nucleotide residues, single-stranded or double-stranded DNA or RNA.

The experimental section below describes a set of libraries of zinc finger genes which code for a set of libraries of zinc finger proteins,  
15 which are used to identify specific zinc fingers which interact with specific polynucleotides. But the invention is more broadly applicable. It is in principle possible to make a set of libraries of any protein which undergoes a specific binding interaction, using that protein as a scaffold to vary specific amino acid residues. It is in principle possible to make a set of  
20 libraries of genes coding for such a set of protein libraries. And it is possible to use such a set of protein libraries to investigate any specific binding interaction, e.g. where the specific binding partner is a polynucleotide or another protein or a different molecule. It may be noted that zinc fingers may be capable of undergoing specific binding  
25 interactions, not only with polynucleotides, but also with other proteins.

It is convenient to control the overlap between libraries of a set of protein libraries by controlling the DNA sequences of the genes which code for the proteins. Thus, to make a library of zinc finger proteins, a library of zinc finger genes is first made. For convenience in relation to  
30 what follows we quote the genetic code which relates the identities of codons to the amino acids which they specify.

- 9 -

		2nd base					
		A C G T					
1st base	A	Lys	Thr	Arg	Ile	3rd base	A C C T
		Asn	Thr	Ser	Ile		
		Lys	Thr	Arg	Met		
		Asn	Thr	Ser	Ile		
	C	Gln	Pro	Arg	Leu		A C G T
		His	Pro	Arg	Leu		
		Gln	Pro	Arg	Leu		
		His	Pro	Arg	Leu		
	G	Glu	Ala	Gly	Val		A C G T
		Asp	Ala	Gly	Val		
		Glu	Ala	Gly	Val		
		Asp	Ala	Gly	Val		
	T	STOP	Ser	STOP	Leu		A C C T
		Tyr	Ser	Cys	Phe		
		STOP	Ser	Trp	Leu		
		Tyr	Ser	Cys	Phe		

Thus for example a codon with multiple degeneracy, e.g.

- 5 ANN comprises 16 different triplets and codes for seven different amino acids namely Lys, Asn, Thr, Arg, Ser, Ile and Met.

- While it is possible in principle to use as few as six libraries of genes to identify a particular amino acid residue, it is in practice convenient to use twelve such libraries in groups of four, wherein libraries 1 to 4
- 10 identify the first nucleotide of a triplet, libraries 5 to 8 identify the second nucleotide of the triplet, and libraries 9 to 12 identify the third nucleotide of the triplet which codes for the amino acid. In this arrangement it is preferable that only one of libraries 1 to 4 (and correspondingly only one of libraries 5 to 8 and only one of libraries 9 to 12) codes for any particular
- 15 amino acid. These considerations give rise to various possible sets of 12 libraries of which one is shown in the following Table 1.

- 10 -

**Table 1**

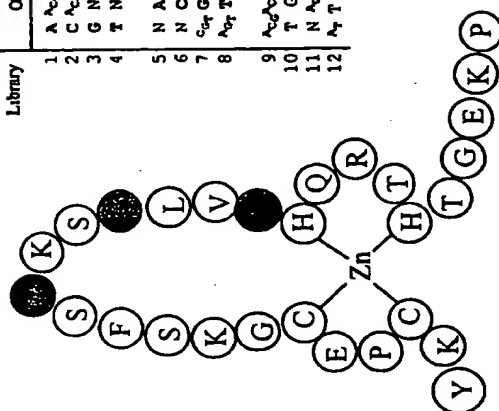
Library	Residue	Codon	Amino Acids Specified
1	$\alpha$	A <sup>A</sup> <sub>C</sub> T N	Lys Asn Thr Ile Met
2	$\alpha$	C <sup>A</sup> <sub>C</sub> G N	Gln His Pro Arg
3	$\alpha$	G N N	Au Asp Ala Gly Val
4	$\alpha$	T N N	Tyr Ser Cys Trp Leu Phe
5	$\alpha$	N A N	Lys Asn Gln His Glu Asp Tyr
6	$\alpha$	N C N	Thr Pro Ala Ser
7	$\alpha$	C <sup>G</sup> <sub>T</sub> G N	Arg Gly Cys Trp
8	$\alpha$	A <sup>C</sup> <sub>T</sub> T N	Ile Met Leu Val Phe
9	$\alpha$	A <sup>C</sup> <sub>G</sub> A <sup>C</sup> <sub>T</sub> G	Lys Thr Met Gln Pro Leu Glu Ala Val
10	$\alpha$	T G G	Trp
11	$\alpha$	N A <sup>G</sup> C	Asn Ser His Arg Asp Gly Tyr Cys
12	$\alpha$	A <sup>T</sup> T C	Ile Phe

5                      Note that any given amino acid appears only once in any set of 4 libraries.

                    Similar randomisation can now be applied to all three positions:  $\alpha$ ,  $\beta$  and  $\gamma$  of zinc finger proteins, to generate libraries 1-36. In libraries 1-12, the randomisation of residue  $\alpha$  is controlled (in these  
 10    libraries, residues  $\beta$  and  $\gamma$  are fully randomised - they are specified by the codon NNN). Similarly, libraries 13-24 control the randomisation of position  $\beta$ , and libraries 25-36 control the randomisation of residue  $\gamma$ ).

- 11 -

Library	$\alpha$	$\beta$	$\gamma$	Library	$\alpha$	$\beta$	$\gamma$
1	A <sub>0</sub> T <sub>0</sub> N	NNN	NNN	25	NNN	NNN	A <sub>0</sub> T <sub>0</sub> N
2	C <sub>0</sub> A <sub>0</sub> G	NNN	NNN	26	NNN	NNN	C <sub>0</sub> A <sub>0</sub> G
3	G <sub>0</sub> N <sub>0</sub> N	NNN	NNN	27	NNN	NNN	G <sub>0</sub> N <sub>0</sub> N
4	T <sub>0</sub> NN	NNN	NNN	28	NNN	NNN	T <sub>0</sub> NN
5	N <sub>0</sub> A <sub>0</sub> N	NNN	NNN	29	NNN	NNN	N <sub>0</sub> A <sub>0</sub> N
6	N <sub>0</sub> C <sub>0</sub> N	NNN	NNN	30	NNN	NNN	N <sub>0</sub> C <sub>0</sub> N
7	C <sub>0</sub> T <sub>0</sub> G	NNN	NNN	31	NNN	NNN	C <sub>0</sub> T <sub>0</sub> G
8	A <sub>0</sub> T <sub>0</sub> T	NNN	NNN	32	NNN	NNN	A <sub>0</sub> T <sub>0</sub> T
9	A <sub>0</sub> C <sub>0</sub> A <sub>0</sub> G	NNN	NNN	33	NNN	NNN	A <sub>0</sub> C <sub>0</sub> A <sub>0</sub> G
10	T <sub>0</sub> G <sub>0</sub> G	NNN	NNN	34	NNN	NNN	T <sub>0</sub> G <sub>0</sub> G
11	N <sub>0</sub> A <sub>0</sub> C	NNN	NNN	35	NNN	NNN	N <sub>0</sub> A <sub>0</sub> C
12	A <sub>0</sub> T <sub>0</sub> C	NNN	NNN	36	NNN	NNN	A <sub>0</sub> T <sub>0</sub> C

Nucleotide sequences of randomised codons  $\alpha$ ,  $\beta$  and  $\gamma$  in libraries 1-36

Randomisation Strategy A

- 12 -

All 36 gene libraries are expressed to generate zinc finger libraries. These zinc finger libraries are then incubated with a polynucleotide of interest, in such a way as to identify one library from each group of four that binds most strongly to the polynucleotide. For example,  
5 each library may be placed in an individual well of a microtitre plate and there incubated with the same trinucleotide.

Consider the controlled randomisation of residue  $\alpha$ . Because in any one group of 4 libraries each amino acid is encoded only once, each amino acid, as residue  $\alpha$ , will occur in only three of the twelve libraries:

10



- 13 -

Library	Lys	Asn	Thr	Ile	Met	Gln	His	Pro	Glu	Asp	Ala	Gly	Val	Tyr	Cys	Trp	Phe	Ser	Arg	Leu
1	✓	✓	✓	✓	✓	-	✓	-	-	-	-	-	-	-	-	-	-	-	-	-
2	-	-	-	-	-	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
5	✓	✓	-	-	-	✓	✓	-	✓	✓	-	-	-	✓	-	-	-	✓	-	-
6	-	-	✓	-	-	-	-	✓	-	-	✓	-	-	-	-	-	-	-	-	-
7	-	-	-	-	-	-	-	-	-	-	-	✓	-	-	✓	-	-	-	✓	-
8	-	-	-	✓	✓	-	-	-	-	-	-	-	✓	-	-	-	✓	-	-	✓
9	✓	-	✓	-	✓	✓	-	✓	✓	-	✓	-	✓	-	-	-	-	-	-	✓
10	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
11	-	✓	-	-	-	-	✓	-	-	✓	-	✓	-	✓	-	-	-	✓	✓	-
12	-	-	-	✓	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

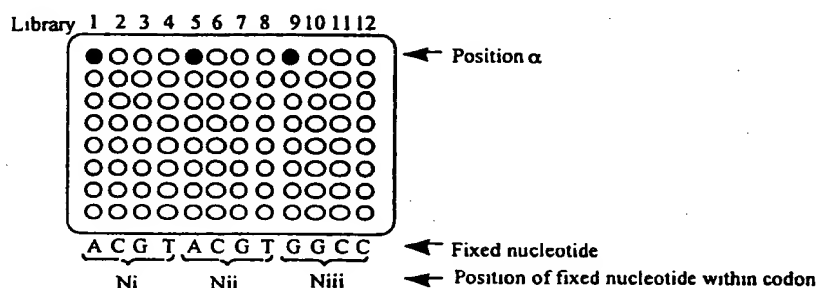
Key:

✓: Specified amino acid is present in this library, at position  $\alpha$ .-: Specified amino acid is not present in this library, at position  $\alpha$ .

- 14 -

Presence / absence of an amino acid at position  $\alpha$  within any given library is a direct result of the controlled randomisation and the genetic code.

- This may now be applied to the assay. Consider that libraries 1-12 only are screened with the trinucleotide ATG and that in order for a zinc finger to bind ATG, residue  $\alpha$  must be Lys (lysine). An assay of libraries 1-12 is performed:



10

- Only libraries 1, 5 and 9 contain lysine as residue  $\alpha$ , therefore only these libraries can emit light. None of the other libraries can emit light, because none of them specify lysine as residue  $\alpha$ . However, this is not the limit of our knowledge. We know the identity of the fixed nucleotide within each library. Moreover, we can read this off directly from the microtitre plate. In this case, the order of fixed nucleotides is AAG.

- Thus, simply from the unique combination of libraries which emit light, we know the genetic code for the amino acid required as residue  $\alpha$ . In this case, the essential fixed nucleotides are AAG, which specifies lysine. We have now linked the genetic code directly to the physical properties of a protein.

- This principle may be applied to all 36 libraries. In so doing, the genetic codes and thus required identities of all three residues  $\alpha$ ,  $\beta$  and  $\gamma$  will be determined:



- 16 -

This is possible, because in libraries 1-12, residues  $\beta$  and  $\gamma$  are fully randomised. Therefore, in each of libraries 1-12 Ser and Arg are present as residues  $\beta$  and  $\gamma$  within the mixture.

Similarly, when controlled randomisation is applied to residue  $\beta$  (libraries 13-24) residues  $\alpha$  and  $\gamma$  are fully randomised and when controlled randomisation is applied to residue  $\gamma$ , residues  $\alpha$ ,  $\beta$  are fully randomised.

By screening the 36 libraries with each of the 64 trinucleotides, an optimum zinc finger will be found for each trinucleotide. Thus the result is therefore the solution of the zinc finger code whereby DNA binding proteins may now be designed at will.

The above strategy generates libraries of genes which when expressed, yield protein libraries in which two positions are fully randomised and one position has controlled randomisation. In practice, this leads to libraries with between 400 (e.g. library 10) and 3600 (eg. library 9) constituent proteins. These numbers are calculated as follows:

$$\begin{aligned}
 \text{Number of library constituents} &= \text{multiplication of number of possibilities at each position of randomisation} \\
 \text{eg. library 1:} &= \text{position } \alpha \times \text{position } \beta \times \text{position } \gamma \\
 &\quad 5 \quad \times \quad 20 \quad \times \quad 20 \\
 &= \underline{\underline{2000 \text{ constituents (proteins)}}}
 \end{aligned}$$

However, these small libraries result from the degeneracy of the genetic code. In practice, the gene libraries which encode the proteins, randomised as above, will be far larger. For example, again consider library 1:

- 17 -

Codon	$\alpha$		$\beta$		$\gamma$
Sequence	A A <sub>C</sub> T N		N N N		N N N
Numbers	1x3x4	x	4x4x4	x	4x4x4 = 49152 constituents (genes)

5

The generation of such libraries should not be problematic technically, since libraries far larger than these exist already (eg. Choo and Klug, 1994, PNAS 91, 11163-7). However, it may it may prove beneficial to reduce the gene library sizes to those of the protein libraries. Potential

10

benefits include:

- greater likelihood of full representation within each library (all constituent proteins encoded);
- even representation of each constituent (an equal amount of each constituent protein within a given library);
- consistent optimum codon usage (to maximise expression).

15

These attributes are desirable because of the degeneracy of the genetic code. Again consider library 1. Within this library, position  $\beta$  is encoded by NNN. When expressed therefore, residue  $\beta$  is 6 times more likely to be serine than it is to be methionine, because serine is encoded six

20

times within NNN for each encoding of methionine.

Such bias within libraries may have an adverse effect on the results of the assay. Any detrimental effect is predicted to be minor - it should occur only if two proteins have similar binding affinities with a given DNA sequence. However, such an eventuality is possible: consider that

two zinc fingers with positions  $\alpha$ =Arg,  $\beta$ =Ser,  $\gamma$ =Lys and  $\alpha$ =Arg,  $\beta$ =Met,  $\gamma$ =Lys bind similarly to a given sequence of DNA, with  $\alpha$ =Arg,  $\beta$ =Met,  $\gamma$ =Lys being the optimally binding zinc finger protein. During the assay, the effective concentration of the protein containing serine at position  $\beta$  would be greater than that of the protein containing methionine. Thus, the serine-

containing protein might give a stronger signal even though it is not the optimum zinc finger for that DNA sequence.

30

- 18 -

It may therefore be preferred to substitute the codon MAX for positions of full randomisation (previously NNN), where MAX is a mixture containing **only** the following codons:

- 5    AAA, AAC, ACC, AGC, ATG, ATT, CAG, CAT, CCG, CGC, CTG, GAA, GAT, GCG, GGC, GTG,  
TAT, TGG, TGC, TTT.

- These codons represents those most favoured by *E. coli* for each amino acid (Nakamura et al., (1997), Nucleic Acids Research, 25,  
10    244-245).

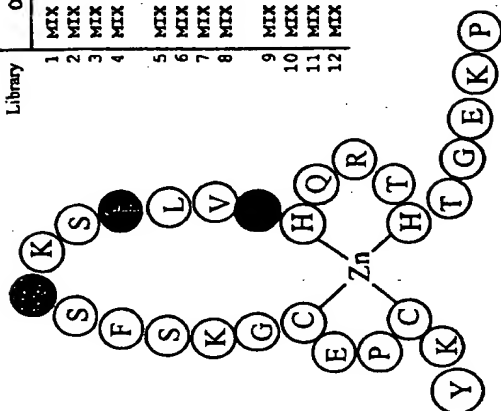
In order to employ these codons in controlled randomisation, a new division of the codons into sets of 12 libraries is required, as outlined in randomisation strategy B.

- 19 -

Library	$\alpha$	$\beta$	$\gamma$	Library	$\alpha$	$\beta$	$\gamma$	Library	$\alpha$	$\beta$	$\gamma$
1	MIX 1	M A X	M A X	13	M A X	MIX 1	M A X	25	M A X	M A X	MIX 1
2	MIX 2	M A X	M A X	14	M A X	MIX 2	M A X	26	M A X	M A X	MIX 2
3	MIX 3	M A X	M A X	15	M A X	MIX 3	M A X	27	M A X	M A X	MIX 3
4	MIX 4	M A X	M A X	16	M A X	MIX 4	M A X	28	M A X	M A X	MIX 4
5	MIX 5	M A X	M A X	17	M A X	MIX 5	M A X	29	M A X	M A X	MIX 5
6	MIX 6	M A X	M A X	18	M A X	MIX 6	M A X	30	M A X	M A X	MIX 6
7	MIX 7	M A X	M A X	19	M A X	MIX 7	M A X	31	M A X	M A X	MIX 7
8	MIX 8	M A X	M A X	20	M A X	MIX 8	M A X	32	M A X	M A X	MIX 8
9	MIX 9	M A X	M A X	21	M A X	MIX 9	M A X	33	M A X	M A X	MIX 9
10	MIX 10	M A X	M A X	22	M A X	MIX 10	M A X	34	M A X	M A X	MIX 10
11	MIX 11	M A X	M A X	23	M A X	MIX 11	M A X	35	M A X	M A X	MIX 11
12	MIX 12	M A X	M A X	24	M A X	MIX 12	M A X	36	M A X	M A X	MIX 12

Nucleotide sequences of randomised codons  $\alpha$ ,  $\beta$  and  $\gamma$  in libraries 1-36

## Randomisation strategy B



- 20 -

where mixes 1 to 12 are as detailed in Table 2:

Mix	1	2	3	4	5	6	7	8	9	10	11	12
Codons	AAA AAC ACC AGC ATG ATT	CAG CAT CCG CGC CTG	GAA GAT GCG GGC GTG	TAT TGC TGG TTT	AAA AAC CAG CAT GAA GAT TAT	ACC CCG GCG GGC	AGC CGC GGC TGC TGG	ATG ATT CTG TTT	AAA GAA AGC CGC GGC TGC	AAC ACC AGC CTG GCG GTG	ATG CAG CCG TAT GCG TGG	ATT CAT GAT TTT
Fixed nucleotide	A	C	G	T	A	C	G	T	A	C	G	T
Position	Ni				Nii				Niii			



- 21 -

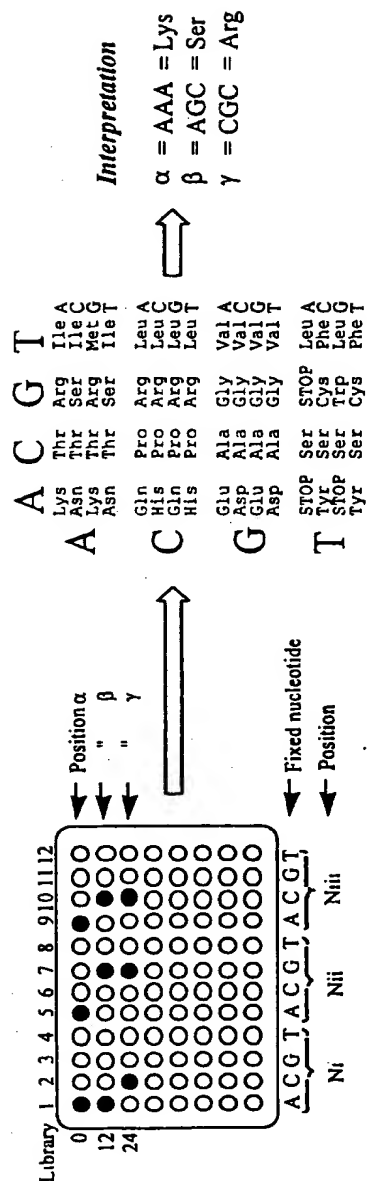
Consider the controlled randomisation of position  $\alpha$  (libraries 1-12). When expressed, position  $\alpha$  will be represented as follows in each library, while positions  $\beta$  and  $\gamma$  are fully randomised.

Library	Lys	Asn	Thr	Ile	Met	Gln	His	Pro	Glu	Asp	Ala	Gly	Val	Tyr	Cys	Trp	Phe	Ser	Arg	Leu
1	✓	✓	✓	✓	✓	-	-	-	-	-	-	-	-	-	-	-	-	✓	-	-
2	-	-	-	-	-	✓	✓	✓	-	-	-	-	-	-	-	-	-	-	✓	-
3	-	-	-	-	-	-	-	-	✓	✓	✓	✓	✓	-	-	✓	-	-	-	-
4	-	-	-	-	-	-	-	-	-	-	-	-	-	✓	✓	✓	✓	-	-	-
5	✓	✓	-	-	-	✓	✓	-	✓	✓	-	-	-	✓	-	-	-	-	-	-
6	-	-	✓	-	-	-	-	✓	-	-	✓	-	-	-	-	-	-	-	✓	-
7	-	-	-	-	-	-	-	-	-	-	-	✓	-	-	✓	✓	-	✓	-	-
8	-	-	-	✓	✓	-	-	-	-	-	-	-	✓	-	-	-	✓	-	-	✓
9	✓	-	-	-	-	-	-	-	✓	-	-	-	-	-	-	-	-	-	-	-
10	-	✓	✓	-	-	-	-	-	-	-	-	✓	-	-	✓	-	-	✓	-	-
11	-	-	-	-	✓	✓	-	✓	-	-	✓	-	✓	-	-	✓	-	-	-	-
12	-	-	-	✓	-	-	✓	-	-	✓	-	-	-	✓	-	-	✓	-	-	-

- 22 -

The changes in controlled randomisation will affect the library numbers which light up and therefore the interpretation of the SPA assay results. However, the principles of controlled randomisation and the mechanism of assay interpretation remain unchanged. Using  
5 randomisation strategy B, the example illustrated above is reiterated:

- 23 -



Note the different fixed nucleotides in libraries 9-12 and that different libraries now light up. The end result:

$\alpha$ =Lys,  $\beta$ =Ser,  $\gamma$ =Arg is the same, however.

- 24 -

Randomisation strategy A is in principle, the easier strategy to implement technically. However, strategy B is preferred. Gene libraries of much smaller size are required. Although construction of these highly-controlled libraries is technically demanding, it is much more likely that the libraries encode all required proteins and moreover that those proteins are encoded in similar proportions, so removing potential difficulties in the SPA library assays.

The above strategies A and B involve designing sets of libraries of genes. It is alternatively possible to devise a set of libraries of proteins. Such a set may consist of

- a) 20 libraries in which each library has one specified amino acid residue at the said first determined position and is randomised at the said one or more other determined positions, and
- b) 20 libraries of corresponding design for each of the said one or more other determined positions.

The method of the invention involves incubating a set of libraries of proteins with a specific binding partner, observing specific binding interactions with certain libraries of the set, and using the observations to identify a protein which interacts with the specific binding partner. Although other assay techniques are possible, this method is preferably performed using scintillation proximity assay (SPA) technology. Briefly, this technology involves providing a support which comprises a scintillant which emits light when subjected to electrons (e.g.  $\beta$  particles) or other forms of radiation resulting from decomposition of a radioisotope. The support may be massive, e.g. the base of each well of a microtitre plate, or may be particulate. One assay reagent is immobilised on the support. Another assay reagent is radiolabelled and is partitioned between two fractions, one bound to the support and the other free in solution. The relative size of the two fractions is arranged to be related to the presence or the concentration of an analyte of interest. The radioisotope is chosen such that reagent bound to the support causes the scintillant in the support

- 25 -

to emit light, while reagent free in solution does not (on account of the short mean free path of the radiation) significantly affect the scintillant substance.

Various assay formats are possible. For example, each library of a set of libraries can be immobilised in an individual well, either  
5 using a scintillant-containing bead (fluomicrosphere) or onto the scintillant base, of a microtitre plate. A specific binding partner of the proteins is radiolabelled and introduced into each well. A specific binding interaction can be investigated in real time. Where several wells emit light, repeated  
10 washing can be used to remove weakly interacting species until the specific binding partner remains bound only in a single well. This ability to identify a single library (as opposed to a small pool of libraries) that bind most strongly to any particular specific binding partner, is a valuable feature, and an advance on assay techniques used previously for similar purposes.

Alternatively, the specific binding partner can be immobilised  
15 in each well of the SPA microtitre plate. Each protein library is radiolabelled and introduced into a different well of the plate for interaction with the specific binding partner. Alternative assay formats, in which neither the protein library nor its specific binding partner, but rather a third reagent is radiolabelled, are well known in the art.

20 Techniques for immobilising protein or other assay reagents on SPA surfaces in forms suitable for taking part in SPA assays, are well known in the art. Development of suitable techniques should not amount to more than the routine optimisation ordinarily required for assays of this kind.

25 Most zinc finger proteins form the DNA recognition module of transcription factors, which serve to switch genes on or off. Already, several examples exist where novel transcription factors have been engineered, by changing their zinc fingers (Choo *et al* (1994), Nature 372, 642-5). Similarly, zinc fingers have been linked to restriction endonuclease  
30 cleavage domains, to generate novel restriction endonucleases (e.g. Kim *et al* (1996), PNAS 93, 1156-60). The application of zinc fingers is almost

- 26 -

limitless - when ever a need arises to link something to a specific sequence of DNA, it can be met with a series of zinc fingers. However, in order to design DNA-binding proteins at will, there must be available one zinc finger for each trinucleotide. This invention provides enabling technology to  
 5 achieve that object.

### Example

The example involves a single protein, comprising three zinc fingers. Controlled randomisation is applied only to the central zinc finger.  
 10 The two outer zinc fingers are present simply to ensure correct registry with the target DNA sequence and to increase overall binding strength (Choo and Klug, (1994) PNAS 01, 11163-67; Berg (1997) Nature Biotech. 15, 323).

The work is divided into four stages: gene synthesis, gene  
 15 expression, SPA assay formats, SPA results and proof of principle (both current and planned).

### **Gene Synthesis:**

A gene was designed and synthesised to encode the protein

20 T G E K P Y K C P E C G K S F S K K S H L V A H Q R T H  
T G E K P Y K C P E C G K S F S K K S H L V A H Q R T H  
 25 T G E K P Y K C P E C G K S F S K K S H L V A H Q R T H

### KEY:

- X linker residues
- X zinc co-ordinating residues
- 30 X DNA-contacting residues ( $\alpha$ ,  $\beta$  and  $\gamma$ ) (positions -1, +3 and +6)

- 27 -

This protein corresponds to three repeats of Berg's consensus zinc finger sequence (Krizek *et al.*, (1991) JACS 113, 4518-23), with DNA-contacting residues from the first zinc finger of transcription factor Sp1 (Berg (1992) PNAS 89, 11109-10; Shi and Berg, (1995) Chem & Biol. 2, 83-89). Each zinc finger sequence is preceded by a *Kruppel*-type linker peptide (Choo and Klug (1993) NAR 21, 3341-6). By analogy to previous precedent (Shi and Berg, 1995), the three repeats of this novel zinc finger peptide are expected to bind to the dsDNA sequence 5'-GGG GGG GGG-3'.

To maximise gene expression, on converting the sequence into DNA, *E. coli* codon preference was employed (Wada *et al.* (1992) NAR20 sup., 2111-8). Wherever possible, first preference codons were used. However, in some instances, second preference codons were also employed. These limited sequence repetition within the gene, necessary to prevent potential intragenic recombination events, which would be deleterious to ensuing experiments. In practice, a maximum repeat length of 8 base pairs was mostly achieved. Use of second preference codons also allowed the incorporation of restriction enzyme sites within the gene. The final gene sequence, restriction sites and codon usage are illustrated in Figure 1.

### Gene Expression

In the current assay format, the zinc finger gene is fused to the glutathione-S-transferase gene in the vector pGEX2TK (Amersham Pharmacia Biotech). Expression of this construct leads to a 36.5 kD protein comprising GST at the amino terminus and the zinc finger protein at the carboxyl terminus. Gene expression is performed in *E. coli* BL21 cells according to manufacturer's instructions. The resulting fusion protein is then purified using glutathione-Sepharose (Amersham Pharmacia Biotech) according to manufacturer's instructions. Use of the pGEX2TK vector allows for the subsequent radiolabelling of the protein if required.

- 28 -

**SPA Assay Format**

The assay format is analogous to that reported for the SPA-based assay NF $\kappa$ B (Amersham Life Science, Proximity News, October 1995). In this format,

- 5 (i) Protein A-derivatised PVT SPA beads bind the Fc region of anti-GST antibody.
- (ii) Anti-GST antibody binds GST of the GST-ZF fusion protein.
- (iii) The GST-ZF fusion protein binds  $^{33}\text{P}$ -labelled oligonucleotide, enabling detection of the protein-DNA interaction.

10

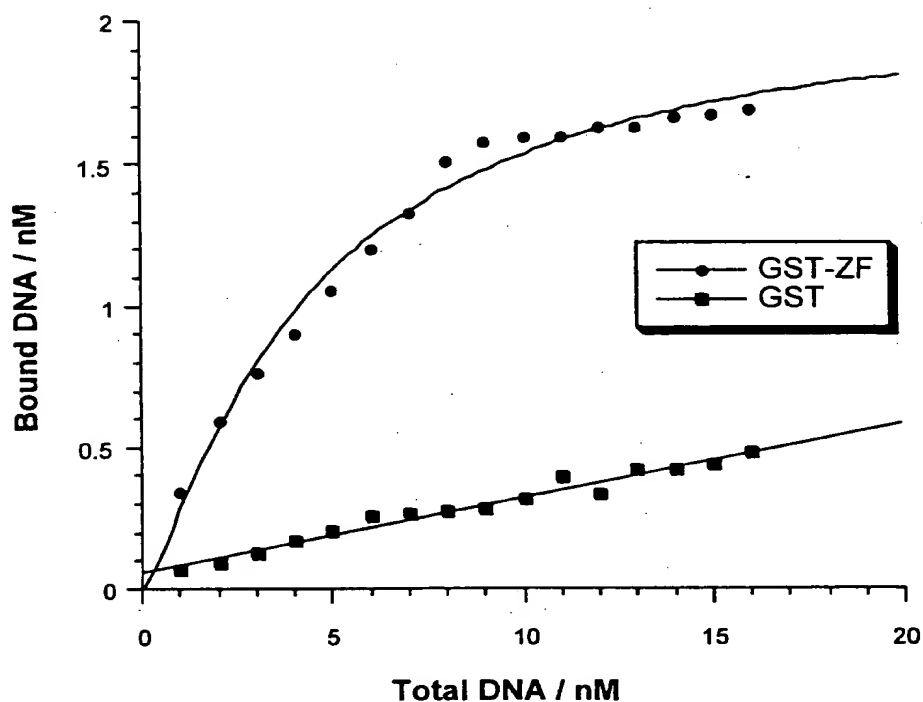
**SPA Results**

- A range of target DNA concentrations was tested in the assay format described above, using either GST-ZF protein, or GST protein as a negative control. Nominal protein concentrations (as estimated by the
- 15 BioRad protein assay with BSA as a reference) were 15nM. Specific activity of the DNA was 247469 cpm/pmol.



- 29 -

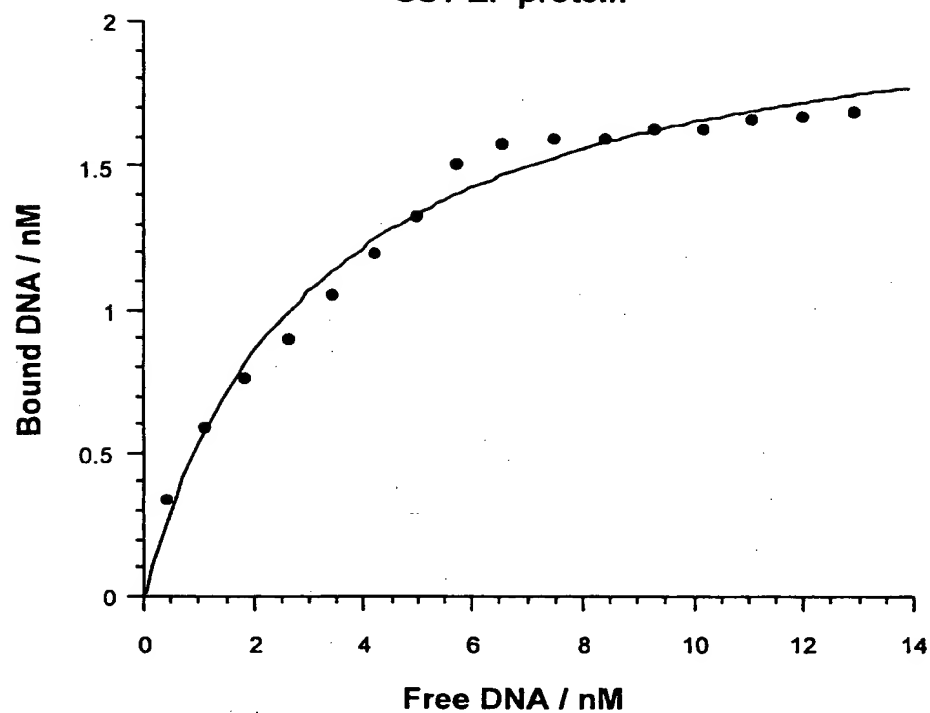
Graph to show [Bound DNA] v. [Total DNA]



To calculate the dissociation constant for the DNA target sequence with the GST-ZF protein, bound DNA is plotted against free DNA. This may be performed either with, or without the GST background being subtracted:

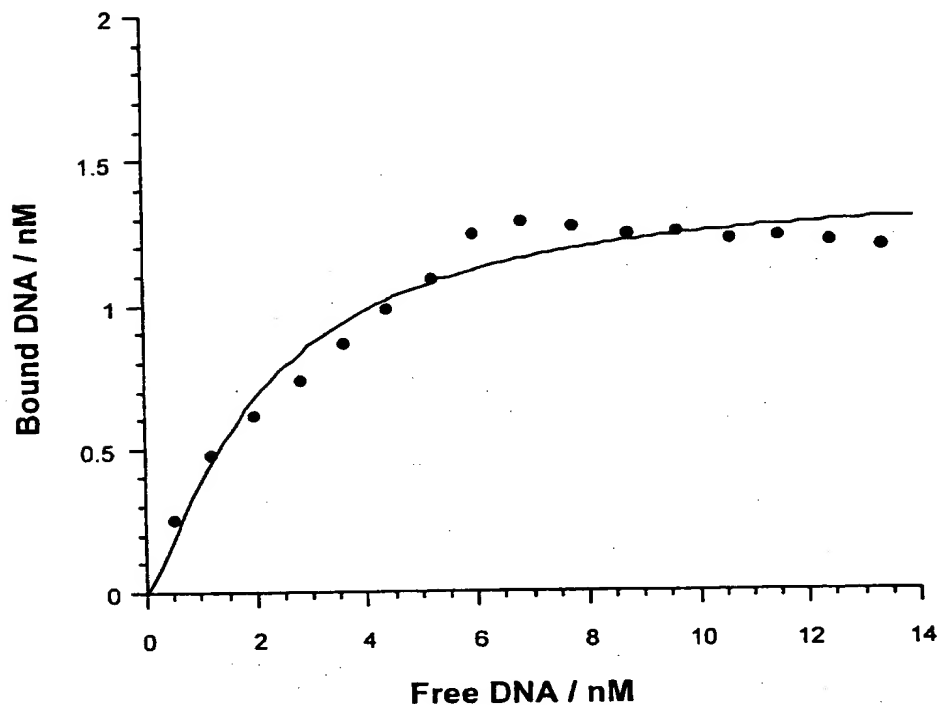
- 30 -

Graph to show [Bound DNA] v. [Free DNA]  
GST-ZF protein



- 31 -

Graph to show [Bound DNA] v. [Free DNA]  
GST-ZF protein, GST background subtracted



#### Calculation of Dissociation Constants

5

Dissociation constants were calculated from the equation

$$\text{Bound ligand} = \frac{B_{\max} * \text{Free ligand}^n}{\text{Free ligand}^n + K_d^n}$$

10 Data from experiments were plotted using KaleidaGraph  
(Synergy software) and the data fitted according to the above equation,  
where:

- 32 -

$m_0$  = Free ligand

$m_1$  =  $B_{\max}$

$m_2$  =  $n$  (Hill coefficient)

$m_3$  =  $K_d$

5 The resulting equations are given below, both for the data without the background subtracted and for data with the background subtracted.

Background not subtracted

$y = m_1 \cdot m_0^m / (m_0^m + m_3)$		
	Value	Error
$m_1$	2.185	0.21632
$m_2$	0.98026	0.14215
$m_3$	3.0801	0.36581
Chisq	0.063505	NA
R	0.98887	NA

10

Background subtracted

$y = m_1 \cdot m_0^m / (m_0^m + m_3)$		
	Value	Error
$m_1$	1.4141	0.098969
$m_2$	1.3027	0.23626
$m_3$	2.6955	0.44464
Chisq	0.076384	NA
R	0.97563	NA

15

Thus, it may be seen that a low nM dissociation constant (approximately 3 nM) between the zinc finger protein and its target DNA sequence has been demonstrated.

- 33 -

CLAIMS

- 5 1. A set of libraries of genes which code for proteins which are capable of specific binding interactions by virtue of amino acid residues at two or more determined positions including a first determined position and one or more other determined positions, which set of libraries consists of:
- 10 a) 6 to 20 libraries in which each library has a triplet that codes for one or several but less than 20 amino acids at the said first determined position, and is randomised at the triplet or triplets coding for the said one or more other determined positions, the arrangement being such that interactions of the proteins coded for by the said 6 to 20 libraries with a specific binding partner identifies a triplet that codes for an amino acid at
- 15 the said first determined position that takes part in the specific binding interaction, and
- b) 6 to 20 libraries of corresponding design for each of the said one or more other determined positions.
- 20 2. The set of libraries of genes as claimed in claim 1, which set of libraries consists of:
- a) 12 libraries in which each library has a triplet that codes for one or several but less than 20 amino acids at the said first determined position, the triplets being as shown in Table 1 or Table 2, and
- 25 b) 12 libraries of corresponding design for each of the said one or more other determined positions.
3. The set of libraries of genes as claimed in claim 1 or claim 2, wherein the genes code for zinc fingers.
4. The set of libraries of genes as claimed in claim 3, which set consists of 36 libraries in three groups of 12 libraries which code for amino
- 30 acids at the -1 and +3 and +6 positions respectively.

- 34 -

5. The set of libraries of genes as claimed in claim 3 or claim 4, wherein each gene codes for a protein comprising 3 zinc fingers.

6. The set of libraries of genes as claimed in claim 5, wherein each gene codes for a protein having the sequence

5  
T G E K P Y K C P E C G K S F S X K S X L V X H Q R T H  
T G E K P Y K C P E C G K S F S X K S X L V X H Q R T H  
10 T G E K P Y K C P E C G K S F S X K S X L V X H Q R T H.

where X is any amino acid

7. A set of libraries of proteins, which proteins are capable of specific binding interactions by virtue of amino acid residues at two or more  
15 determined positions including a first determined position and one or more other determined positions, which set of libraries consists of:

a) 6 to 20 libraries in which each library has one or several but less than 20 amino acid residues at the said first determined position and is randomised at the said one or more other determined positions, the  
20 arrangement being such that interaction of the 6 to 20 libraries with a specific binding partner identifies an amino acid residue at the said first determined position that takes part in the specific binding interaction, and  
b) 6 to 20 libraries of corresponding design for each of the said one or more other determined positions.

25 8. The set of libraries of proteins as claimed in claim 7, which set of libraries consists of

a) 20 libraries in which each library has one specified amino acid residue at the said first determined position and is randomised at the said one or more other determined positions, and  
30 b) 20 libraries of corresponding design for each of the said one or more other determined positions.

- 35 -

9. The set of libraries of proteins as claimed in claim 7 or claim 8, wherein the proteins are zinc fingers.

10. The set of libraries of proteins as claimed in claim 7, which set consists of 60 libraries in three groups of 20 libraries with specified amino acids at the -1 and +3 and +6 positions respectively.

11. The set of libraries of proteins as claimed in claim 9 or claim 10, wherein each protein comprises three zinc fingers.

12. The set of libraries of proteins as claimed in claim 11, wherein each protein as the sequence

10

T G E K P Y K C P E C G K S F S X K S X L V X H Q R T H

T G E K P Y K C P E C G K S F S X K S X L V X H Q R T H

15

T G E K P Y K C P E C G K S F S X K S X L V X H Q R T H.

where X is any amino acid

13. The set of libraries of proteins as claimed in any one of claims 7 to 12, which set results from expression of the set of libraries of genes as claimed in any one of claims 1 to 6.

14. A set of libraries of genes which code for the set of libraries of proteins defined in any one of claims 7 to 12.

15. A method of identifying a protein which interacts with a specific binding partner, which method comprises providing a set of libraries of proteins as defined in any one of claims 7 to 13, incubating the specific binding partner with each library of the set, observing specific binding interactions with certain libraries of the set, and using the observations to identify a protein which interacts with the specific binding partner.

16. The method as claimed in claim 15, wherein the specific binding partner is a polynucleotide.

- 36 -

17. The method as claimed in claim 15, wherein the specific binding interactions are observed by scintillation proximity assay.

18. The method as claimed in claim 17, wherein the sets of libraries of proteins are immobilised on scintillation proximity assay surfaces and the specific binding partner is radiolabelled.

19. The method of claim 17 or claim 18, wherein after incubation the scintillation proximity assay surfaces are washed to distinguish stronger specific binding interactions from weaker ones.

20. A protein having the sequence

10

T G E K P Y K C P E C G K S F S K K S H L V A H Q R T H

T G E K P Y K C P E C G K S F S K K S H L V A H Q R T H

15

T G E K P Y K C P E C G K S F S K K S H L V A H Q R T H.

21. A gene which codes for the protein of claim 20.



- 37 -

**ABSTRACT**

5

**GENE AND PROTEIN LIBRARIES AND METHODS**

A set of libraries of proteins, which proteins are capable of specific binding interactions by virtue of amino acid residues at two or more determined positions including a first determined position and one or more other determined positions, which set of libraries consists of:

- 10 a) 6 to 20 libraries in which each library has one or several but less than 20 amino acid residues at the said first determined position and is randomised at the said one or more other determined positions, the arrangement being such that interaction of the 6 to 20 libraries with a
- 15 specific binding partner identifies an amino acid residue at the said first determined position that takes part in the specific binding interaction, and
- b) 6 to 20 libraries of corresponding design for each of the said one or more other determined positions.

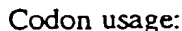
20 A set of libraries of genes which code for the proteins.

A method of identifying a protein which interacts with a specific binding partner, which method comprises incubating the protein with each library of the set of libraries of proteins, observing specific binding interactions with certain libraries of the set, and using the observations to identify a protein which interacts with the specific binding

25 partner.

**This Page Blank (uspto)**

Restriction map:



**86 codons**

MW : 9598 Dalton

TTT phe F	1	TCT ser S	3	TAT tyr Y	2	TGT cys C	1
TTC phe F	2	TCC ser S	3	TAC tyr Y	1	TGC cys C	5
TTA leu L	-	TCA ser S	-	TAA och Z	1	TGA opa Z	-
TTG leu L	-	TCG ser S	1	TAG amb Z	-	TGG trp W	-
CTT leu L	-	CCT pro P	-	CAT his H	6	CGT arg R	1
CTC leu L	-	CCC pro P	-	CAC his H	3	CGC arg R	2
CTA leu L	-	CCA pro P	-	CAA gln Q	1	CGA arg R	-
CTG leu L	3	CCG pro P	6	CAG gln Q	2	CGG arg R	-
ATT ile I	-	ACT thr T	-	AAT asn N	-	AGT ser S	-
ATC ile I	-	ACC thr T	4	AAC asn N	-	AGC ser S	3
ATA ile I	-	ACA thr T	-	AAA lys K	12	AGA arg R	-
ATG met M	-	ACG thr T	2	AAG lys K	3	AGG arg R	-
GTT val V	1	GCT ala A	-	GAT asp D	-	GGT gly G	3
GTC val V	-	GCC ala A	1	GAC asp D	-	GGC gly G	3
GTA val V	-	GCA ala A	1	GAA glu E	4	GGA gly G	-

*This Page Blank (uspto)*